

OPENCUDA+MPI

A FRAMEWORK FOR HETEROGENEOUS GP-GPU DISTRIBUTED COMPUTING

Kenny Ballou
Nilab Mohammad Mousa
College of Engineering – Department of Computer Science
Boise State University
Alark Joshi, Ph.D

Abstract

The introduction and rise of General Purpose Graphics Computing has significantly impacted parallel and high-performance computing. It has introduced challenges when it comes to distributed computing with GPUs. Current solutions target specifics: specific hardware, specific network topology, a specific level of processing. Those restrictions on GPU computing limit scientists and researchers in various ways. The goal of OpenCUDA+MPI project is to develop a framework that allows researchers and scientists to write a general algorithm without the overhead of worrying about the specifics of the hardware and the cluster it will run against while taking full advantage of parallel and distributed computing on GPUs. As work toward the solution continues to progress, we have proven the need for the framework and expect to find the framework enables scientists to focus on their research.

Keywords: Parallel Computing, Distributed Computing, General Purpose Graphics Processing, High-Performance Computing, Scientific Computing

1 Topic

Increasingly, **Graphics Processing Units (GPUs)** are being used for general purpose computation (**General Purpose GPU (GP-GPU)**). They have significantly altered the way high-performance computing tasks can be performed today. To accelerate general-purpose scientific applications, computationally intensive portions of the application are passed to **GPUs**. Once the **GPU** has completed its task it sends the result back to the **Central Processing Unit (CPU)** where the application code is running. This process can make applications run noticeably faster.

Scientific applications require a large number of floating point number operations, **CPUs** are not sufficient to carry out such computationally intensive tasks. **CPUs** are responsible for prioritization and execution of every instruction. Generally, a processor is described by the number of execution cores it owns. Modern **CPUs** have eight cores while **GPUs** have hundreds or more cores. More cores grant **GPUs** the ability to perform more tasks at the same time.

In computer architecture there are two main ways of processing: serial and parallel. **CPUs** consist of a small number of cores (microprocessors) that are best at processing serial data. On the other hand, **GPUs** consist of thousands of cores that are designed for processing parallel data. Given a program, we can run the parallel portions of the code on **GPUs** while the serial portions run on the **CPU**. The programmable **GPU** has evolved into a highly parallel, multithreaded, many core processor with tremendous computational power. **Compute Unified Device Architecture (CUDA)** is an architecture for utilizing and distributing computational tasks onto a computer's **GPUs**. As a parallel computing platform and programming model, **CUDA** utilizes the power of **GPU** to achieve dramatic increases in computing performance [6]. This fact is well illustrated in fig. 1.

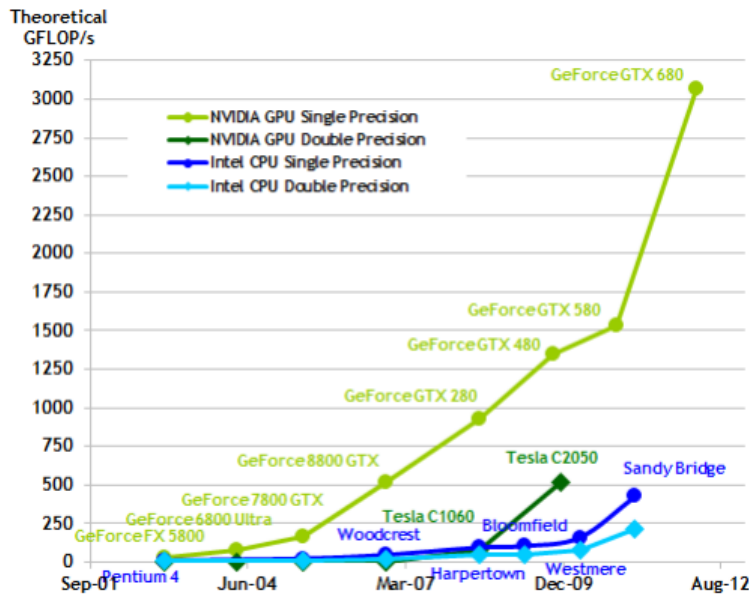


Figure 1: Floating-Point Operations per Second for the **CPU** and **GPU** [6]

Today, more than one million **CUDA-enabled GPUs** are used by software developers, scientists and researchers in order to obtain a large performance gain in wide-ranging applications [6]. A major challenge remains in being able to seamlessly integrate multiple workstations to further parallelize the computational tasks [14] [16]. Current approaches provide the ability to parallelize tasks, but they are less focused on optimally utilizing the varied capabilities of heterogeneous graphics cards in a cluster of workstations (across many computer nodes).

We propose to create a framework for using both **Message Passing Interface (MPI)** and

CUDA on a cluster of computers to dynamically and easily assign computationally intensive tasks to the machines participating in the cluster. **MPI** is a popularly used standardized interface for communication among a group of computers. It facilitates the passing of messages between the nodes in the cluster.

Our framework will be easy to use on a heterogeneous cluster of computers, each containing a **CUDA**-capable **GPU**. The framework will optimize the scheduling of low-level computational tasks based on the capabilities of the varied **GPUs** in the cluster. The framework will provide the system administrator with the ability to easily configure it to allow optimal use of the individual **GPUs**. Overall, the goal is to create a framework that is simple and easy to use, facilitating the combination of two powerful means of computing to further increase the throughput and overall computing power of a cluster.

We plan to evaluate the efficacy of our framework on the problem of vessel extraction. Currently, we use a single **GPU** to extract vascular structures from a CT angiography scan which is computationally intensive [13]. Using the framework to accelerate the computational task will provide us with key insights on its usability.

2 Significance

CUDA+MPI will serve to better utilization of existing computing resources. In other words, it will give us high compute power at low cost. Many personal computers have a **CUDA** graphics card. It is more cost effective to build a cluster of multiple computers than to purchase a supercomputer to perform high performance computing. The reason for the low costs stem from the fact that clusters are built of components which are sold in high volumes [13]. CUDA+MPI will have the advantage of handling heterogeneous distributed computing. In other words, CUDA+MPI avoids hardware limitations. All the computers participating in the cluster do not need to have the same model of **CUDA**-enabled **GPU** cards. This will also give us the freedom to easily add new nodes to the cluster as needed.

Since CUDA+MPI will operate in a heterogeneous cluster it will manage the job requests from the user to be run on a heterogeneous cluster of **GPU**-capable computers. In order to decide job priority, amount of work and availability of devices will be considered. In addition to improved performance, the proposed framework will grant researchers and scientists the assurance that the heterogeneous cluster will be optimized for computing without requiring extra effort and thought.

3 Related Works

Current approaches to accelerating research computations and engineering applications include parallel computing and distributed computing. A parallel system is when several processing elements work simultaneously to solve a problem faster. The primary consideration is elapsed time, not throughput or sharing of resources at remote locations. A distributed system is a collection of independent computers that appears to its users as a single coherent system. In distributed computing the data is split into smaller parts and subsequently spread across the system. The result is then collected and outputted. Summary of current solutions are listed in table 1 below followed by detailed description of every entry.

Project	Description
Hadoop	Distributed computation framework [7] [4] [22]
BOINC	Volunteer and grid computing project distributed [10]
TORQUE	Job scheduler for distributed computing [9]
GPUDirect	GPU-to-GPU framework for parallel computing [3]
MPI	Message-passing system for parallel computations [24] [19]
PVM	Distributed environment and message passing system [5]

Table 1: Summary of Current Approaches and Projects

3.1 Hadoop

Hadoop provides a distributed file system and a framework for the analysis and transformation of very large data sets using the MapReduce paradigm [4] [17] [12] [11]. An important characteristic of Hadoop is the partitioning of data and computation across many (thousands) of hosts (nodes), and executing application computations in parallel close to their data [22]. The MapReduce [20] [8] paradigm which Hadoop implements is characterized by dividing the application into many small fragments of work, each of which may be executed or re-executed on any node in the cluster. The worker node processes the smaller problem, and passes the answer back to its master node. The answers to the subproblems are then combined to form the output. Hadoop is popular model for distributed computations, however, it does not integrate with CUDA-enabled GPUs.

3.2 BOINC

The Berkeley Open Infrastructure for Network Computing (BOINC) is an open source middleware system for volunteer and grid computing. The general public can contribute to today's computing power by donating their personal computer's disk space and some percentage of CPU and GPU processing. In other words, BOINC is software that can use the unused CPU and GPU cycles on a computer to do scientific computing [10]. The BOINC project allows for distributed computing using hundreds of millions of personal computers and game consoles belonging to the general public. This paradigm enables previously infeasible research and scientific super-computing. The framework requires that individuals are connected to the Internet and is supported by various operating systems, including Microsoft Windows, Mac OS X and various Unix-like systems. Although BOINC is another step in grid computing, it does not seem to be relevant to our research.

3.3 TORQUE

Terascale Open-Source Resource and Queue Manager (TORQUE) is a job scheduler, a computer application that controls program execution. It enables users to control computational tasks over distributed compute nodes, as well as scheduling and administration on a cluster [9]. **TORQUE** has the ability to handle larger clusters with tens of thousands of nodes. It also can handle larger jobs that span hundreds of thousands of processors due to advanced **GP-GPU** scheduling. However, **TORQUE** does not allow for heterogeneous hardware in the cluster, making it inferior compared to **CUDA+MPI** framework.

3.4 GPUDirect

Currently **GPU** based clusters are becoming more popular. **GPUDirect** is a **GPU** based clustering using Infiniband, an architecture specification that defines a connection between processor nodes. **GPUDirect** allows for multiple **GPUs** to communicate with each other directly by giving **GPUs** ability to initiate and manage memory transfer [3]. Prior to **GPUDirect**, **GPU-to-GPU** communication involved the **CPU**. **GPUDirect** performance gain exceed **CPU** solutions. Direct **GPU-to-GPU** communication is less computationally expensive since less messages will be sent and received via the host server. An issue with **GPUDirect** is that it requires specific hardware. The proposed framework will serve as software solution to avoid hardware limitations.

3.5 MPI

MPI is the dominant message passing programming paradigm for clusters [1] [2]. **MPI** is a standardized and portable message-passing system that functions on a wide variety of parallel computers. Although the standard **MPI** defines the syntax and semantics of a core of library routines in the **C** programming language, it is a language-independent communications protocol used to program parallel computers. **CUDA+MPI** builds upon this model due to the fact that **MPI** is the dominant model used in high-performance computing [23]. **MPI** has been implemented for almost every distributed memory architecture and is optimized for the hardware on which it runs.

3.6 PVM

Parallel Virtual Machine (PVM) is a software tool for parallel networking of computers. It is designed to allow a network of heterogeneous Unix and/or Windows machines to be used as a single distributed parallel processor. Thus, large computational problems can be solved more cost effectively by using the aggregate power and memory of many computers. **PVM** enables users to exploit their existing computer hardware to solve much larger problems at less additional cost [5]. **PVM** was a step towards modern trends in distributed processing and grid computing but has, since the mid-1990s, largely been supplanted by the much more successful **MPI** standard for message passing on parallel machines.

3.7 CUDA+MPI

The proposed **CUDA+MPI** framework will take advantage of both **MPI** and **CUDA** to perform computations on **GPU** cards of computers participating in the cluster. **CUDA+MPI**'s goal is to allow a collection of heterogeneous computers each containing a **CUDA**-capable **GPU** to be used as a coherent and flexible concurrent computational resource.

4 Methodology

Our framework will be developed using a number of tools and will also build upon a few already developed and mature software libraries. Overall, we plan to use the Python programming language, MPI for cross process communication, and CUDA for graphics computing.

4.1 Python

We will use Python programming language because of its ease of development and plethora of existing libraries such as MPI. Python also adds some advantages when it comes to usability when needing more performance. For example, if we discover a need for certain aspects of the project to be tuned or otherwise run faster, we can easily switch to C or C++ and write sections of the program in a (more performant) native language.

4.2 MPI

Currently our cluster consists of 16 computers provided to our research lab by the computer science department of Boise State University. We will use MPI because it has established itself as the de-facto interface for cross process communication [24]. To allow for interfacing with Python [19], we will use mpi4py because of its maturity and implementation completeness.

4.3 CUDA

Further, we will be using CUDA because of existing knowledge of the framework and also because it is a well established framework for GPU computing.

4.4 CUDA+MPI

Our framework’s goal is to be able to manage the job queue to be run on GPUs for a cluster of GPU-capable computers. As more implementation work is done one or more of the parameters listed below will be considered to decide which particular job to run.

- Job priority
- Compute resource availability
- Execution time allocated to user
- Number of simultaneous jobs allowed for a node
- Estimated execution time
- Elapsed execution time
- Availability of devices

4.5 Testing

To evaluate the efficacy and accuracy of our framework several tests will be executed. As previously mentioned one of the first algorithms to be tested will be the problem of vessel extraction. Accurately extracting vascular structures from a Computerized Tomographic angiography—also called CT angiography scans—is important for creating oncologic surgery planning tools as well as medical visualization applications [13]. Currently, we use a single GPU to extract vascular structures from a CT angiography scan, which is computationally intensive. The following test programs will be developed as time allows:

- N-Body Simulation
- Prime Number Searching

Every test program will be evaluated in three categories: **CPU**-only, **CUDA**-only, and **CUDA+MPI**. Using the framework to accelerate the computational tasks will provide us with key insights on its usability.

5 Results

5.1 Vector Summation

Our first developed test program was a 10^9 element wise vector summation problem. This is a simple and, as we will see, a bad example of using a cluster to speed up the computation of a problem. Although we did see a increase in performance, the cost of **Input/ Output (I/O)** far out weighs the benefit. Specifically, to do the computation on one machine (one **CPU**), it took a **wall time** of 254 seconds (about 4 minutes) and a **CPU time** of 13.7 seconds. Further, to do the computation using a single **GPU** took about 4172 seconds (**wall time**) or about 70 minutes, 13.83 seconds (**CPU time**) while the computing the summation of the cluster took about 3177 seconds (**wall time**) or about 50 minutes, 10.51 seconds (**CPU time**). Our **CPU** only implementation took the least amount of elapsed time, it took the second longest **CPU time**. Our **CUDA** only implementation took the longest in both **wall time** and **CPU time** and our cluster implementation was shortest **CPU time** but seconds longest elapsed time. The benefit of saving an upwards of 3.3 seconds is not worth the extra incurred cost of 2923 seconds. Not so surprisingly though, running the vector summation problem over the cluster *without* using **CUDA** does increase the overall elapsed time of the problem; namely, it took 226 seconds **wall time** (currently the correct **CPU time** is unable to be collected). Further, increasing the number of nodes part of the program’s pool, decreases the **wall time** for each node.

Method	Time (s)	Total Time (s)
CPU Only	13.7	254.13
CUDA (Single Node)	13.83	4172
MPI + CUDA (7 nodes)	10.51	(average) 3177
MPI (7 nodes)		(average) 226
MPI (16 slots)		(average) 149

Table 2: Computational Timing Comparison of 10^9 element wise vector summation

5.2 N-Body Problem

We have several sizes of the N-Body problem that we tested with: 2,001 particles, 20,000 particles, 200,000 particles, 2,000,000, and 20,000,000 particles.

The computational times are for only one time step. The method for computing the gravitational potential is an adaptation of the **Particle-Particle, Particle-Mesh (P3M)** method. The benefit of using this method is we are able to nicely distribute the problem over the **cluster** and / or over a **GPU** (because of memory limitations) while maintaining a respectable accuracy for close body interactions. Further, if a grid contains more bodies than a specified threshold (in our case 200,000), we can further sub-divide the grid to improve performance and maintain accuracy still.

There are other algorithms for computing N-Body problems on the **CPU** that are quite efficient, notably, the Barnes-Hutt Tree algorithm; however, using it would distort and confound the comparisons between **CPU**, **GPU**, and **cluster** implementations; foregoing to mention the complexities of implementing such an algorithm on **GPUs** and over a **cluster**.

5.2.1 N-Body — CPU

In **CPU** tests, we were only able to complete several of the problem sizes; the larger sizes are infeasible. Notably, the smaller sizes were computed in a relatively respectable amount of time. While the bigger sizes were time consuming, still being computed or not even attempted.

Size	User (seconds)	Sys (seconds)	Real (seconds)
2001	10.65	0.00	12.03
20000	861.46	0.00	861.72
200000	109306	18.05	109364
2 million
20 million	N/A	N/A	N/A

Table 3: **CPU** N-Body simulation using particle-particle method

5.2.2 N-Body — GPU

As noted above, the **GPU** (**CUDA**) implementation uses the same computational method (**P3M**). Using the **GPU**, we were able to compute the 20,000,000 size problem and may be able to compute larger sets within a *reasonable* amount of time.

Size	User (seconds)	Sys (seconds)	Real (seconds)
2001	0.68	0.48	01.25
20000	3.41	0.55	04.06
200000	31.06	1.11	32.28
2 million	347	11.93	361
20 million	115.47	120.65	13927

Table 4: Single **GPU** (**CUDA**) N-Body simulation using **P3M** method

A Reference

- [1] MPI: A Message Passing Interface Standard. <http://www.mpi-forum.org/>, June 1995.
- [2] MPI-2: Extensions to the Message Passing Interface. <http://www.mpiforum.org/>, July 1997.
- [3] GPU-Direct Accelerating GPU Cluster Performance. <http://www.youtube.com/watch?v=o5Ir93SA8ZE>, May 2010.
- [4] Apache Hadoop. <http://hadoop.apache.org/>, April 2013.
- [5] Computer Science and Division - PVM: Parallel Virtual Machine. <http://www.csm.ornl.gov/pvm/>, April 2013.
- [6] CUDA C Programming Guide. http://docs.nvidia.com/cuda/pdf/CUDA_C_Programming_Guide.pdf, May 2013.
- [7] Hadoop Project Description. <http://wiki.apache.org/hadoop/ProjectDescription>, April 2013.
- [8] MapReduce. <http://wiki.apache.org/hadoop/MapReduce>, April 2013.
- [9] TORQUE Resource Manager-Adaptive Computing. <http://www.adaptivecomputing.com/products/open-source/torque/>, May 2013.
- [10] David P Anderson. BOINC: A System for Public-Resource Computing and Storage. In *Grid Computing, 2004. Proceedings. Fifth IEEE/ACM International Workshop on*, pages 4–10. IEEE, 2004.
- [11] Dean, Jeffrey and Ghemawat, Sanjay. MapReduce: Simplified Data Processing on Large Clusters, OSDI’04: Sixth Symposium on Operating System Design and Implementation, San Francisco, CA, December, 2004. *S. Dill, R. Kumar, K. McCurley, S. Rajagopalan, D. Sivakumar, ad A. Tomkins, Self-Similarity in the Web, Proc VLDB*, 2004.
- [12] Dias, Jonas and Aveleda, Albino. HPC Environment Management: New Challenges in the Petaflop Era. In *High Performance Computing for Computational Science—VECPAR 2010*, pages 293–305. Springer, 2011.
- [13] Erdt, Marius and Raspe, Matthias and Suehling, Michael. Automatic hepatic vessel segmentation using graphics hardware. In *Medical Imaging and Augmented Reality*, pages 403–412. Springer, 2008.
- [14] Hadri, Bilel and Fahey, Mark and Jones, Nick. Identifying Software Usage at HPC Centers with the Automatic Library Tracking Database. In *Proceedings of the 2010 TeraGrid Conference*, page 8. ACM, 2010.
- [15] Hindman, Benjamin and Konwinski, Andy and Zaharia, Matei and Ghodsi, Ali and Joseph, Anthony D and Katz, Randy and Shenker, Scott and Stoica, Ion. Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center. In *Proceedings of the 8th USENIX Conference on Networked Systems Design and Implementation*, pages 22–22. USENIX Association, 2011.
- [16] Hindman, Benjamin and Konwinski, Andy and Zaharia, Matei and Stoica, Ion. A Common Substrate for Cluster Computing. In *Workshop on Hot Topics in Cloud Computing (HotCloud)*, volume 2009, 2009.
- [17] J. Venner. Pro Hadoop. <http://www.apress.com/>, June 2009.
- [18] Johnson, C.R. Biomedical Visual Computing: Case Studies and Challenges. *Computing in Science & Engineering*, 14(1):12–21, 2012.
- [19] Lisandro Dalcin. MPI for Python. <http://mpi4py.scipy.org/docs/usrman/index.html>, January 2012.

- [20] Luo, Yuan and Guo, Zhenhua and Sun, Yiming and Plale, Beth and Qiu, Judy and Li, Wilfred W. A Hierarchical Framework for Cross-Domain MapReduce Execution. In *Proceedings of the Second International Workshop on Emerging Computational Methods for the Life Sciences*, pages 15–22. ACM, 2011.
- [21] Parker, Christopher and Suleman, Hussein. A Lightweight Web Interface to Grid Scheduling Systems. In *Proceedings of the 2008 Annual Research Conference of the South African Institute of Computer Scientists and Information Technologists on IT Research in Developing Countries: Riding the Wave of Technology*, pages 180–187. ACM, 2008.
- [22] Shvachko, Konstantin V. Apache Hadoop: The Scalability Update. *login: The Magazine of USENIX*, 36:7–13, 2011.
- [23] Sur, Sayantan and Koop, Matthew J and Panda, Dhabaleswar K. High-performance and scalable MPI over InfiniBand with reduced memory usage: an in-depth performance analysis. In *Proceedings of the 2006 ACM/IEEE conference on Supercomputing*, page 105. ACM, 2006.
- [24] Wes. A Comprehensive MPI Tutorial Reference. <http://www.mpitutorial.com/>, 2012.

Glossary

CPU time time (typically in seconds) spent executing the process. 9

BOINC Berkeley Open Infrastructure for Network Computing. 5

cluster A collection of computers networked together, typically with the goal of combining computational power. 9

CPU Central Processing Unit. 2, 5, 6, 8–10

CUDA Compute Unified Device Architecture. 2–10

GP-GPU General Purpose GPU. 2, 6

GPU Graphics Processing Unit. 2–7, 9, 10

I/O Input/ Output. 9

MPI Message Passing Interface. 2, 3, 6, 7

node a computer, member of a cluster. 9

P3M Particle-Particle, Particle-Mesh. 9, 10

PVM Parallel Virtual Machine. 6

slot process space on a node, typically equal to the number of logical CPU cores on the node. 9

TORQUE Terascale Open-Source Resource and QUEue Manager. 6

wall time all elapsed time (typically in seconds) that passes from start to finish of a program. 9